

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開2002-182685

(P2002-182685A)

(43)公開日 平成14年6月26日(2002.6.26)

(51)Int.Cl. <sup>7</sup>	識別記号	F I	テームコード*(参考)
G 1 0 L 15/18		G 0 6 T 7/00	P 5 D 0 1 5
G 0 6 T 7/00			3 0 0 F 5 L 0 9 6
	3 0 0	G 1 0 L 3/00	5 3 7 Z
G 1 0 L 15/00			5 5 1 H
15/24			5 7 1 Q
審査請求 未請求 請求項の数12 O L (全 10 頁)			

(21)出願番号 特願2000-376911(P2000-376911)

(22)出願日 平成12年12月12日(2000.12.12)

(71)出願人 000002185

ソニー株式会社

東京都品川区北品川6丁目7番35号

(72)発明者 中塚 洪長

東京都品川区北品川6丁目7番35号 ソニ

ー株式会社内

(74)代理人 100082131

弁理士 稲本 義雄

Fターム(参考) 5D015 FF00 HH04

5L096 BA05 FA23 GA41 HA04 JA11

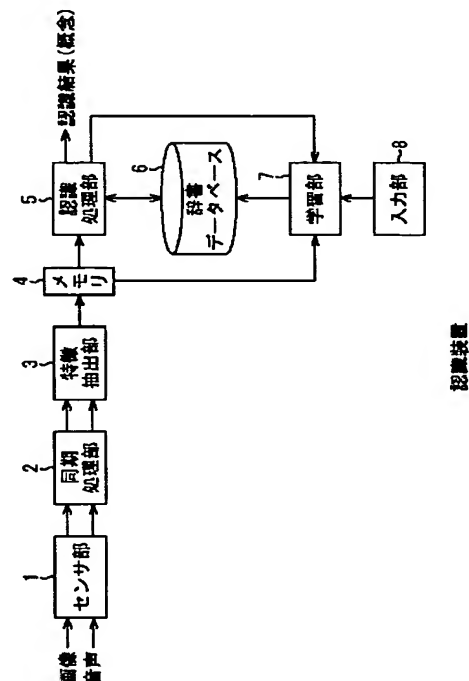
KA04

(54)【発明の名称】 認識装置および認識方法、学習装置および学習方法、並びに記録媒体

(57)【要約】

【課題】 認識性能を向上させる。

【解決手段】 同期処理部2は、入力された画像と音声とを同期させ、特徴抽出部3は、その同期された画像と音声それぞれから、特徴量を抽出して、その画像と音声の特徴量を合成した合成特徴量を得る。学習部7は、その合成特徴量に基づいて学習を行い、同一概念を表す画像および音声に対応するモデルを生成し、そのモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する。一方、認識処理部5は、合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声が表示概念を認識する。



## 【特許請求の範囲】

【請求項 1】 同一概念を表す画像および音声のモデルと、その概念を表す概念情報とを対応付けた辞書を記憶する記憶手段と、

入力された画像と音声とを同期させる同期手段と、同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出手段と、

前記抽出手段において出力される合成特徴量と、前記辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声とが表す概念を認識する認識手段とを備えることを特徴とする認識装置。

【請求項 2】 前記同期手段は、入力された画像の区間である画像区間と、入力された音声の区間である音声区間を検出し、前記画像区間と音声区間それぞれを正規化することにより、入力された画像と音声とを同期させることを特徴とする請求項 1 に記載の認識装置。

【請求項 3】 前記同期手段は、さらに、正規化された前記画像区間と音声区間それぞれの始点と終点とを一致させることにより、入力された画像と音声とを同期させることを特徴とする請求項 2 に記載の認識装置。

【請求項 4】 前記同期手段は、さらに、正規化された前記画像区間と音声区間それぞれにおける画像のフレームと音声のフレームとを対応させることにより、入力された画像と音声とを同期させることを特徴とする請求項 2 に記載の認識装置。

【請求項 5】 同一概念を表す画像および音声のモデルと、その概念を表す概念情報とを対応付けた辞書を参照して認識処理を行う認識方法において、入力された画像と音声とを同期させる同期ステップと、同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、前記抽出ステップにおいて出力される合成特徴量と、前記辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声とが表す概念を認識する認識ステップとを備えることを特徴とする認識方法。

【請求項 6】 同一概念を表す画像および音声のモデルと、その概念を表す概念情報とを対応付けた辞書を参照して行う認識処理を、コンピュータに行わせるプログラムが記録されている記録媒体において、入力された画像と音声とを同期させる同期ステップと、同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、

前記抽出ステップにおいて出力される合成特徴量と、前記辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声とが表す概念を認識する認識ステップとを備えるプログラムが記録されていること

を特徴とする記録媒体。

【請求項 7】 入力された画像と音声とを同期させる同期手段と、

同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出手段と、

前記抽出手段において出力される合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習手段とを備えることを特徴とする学習装置。

【請求項 8】 前記同期手段は、入力された画像の区間である画像区間と、入力された音声の区間である音声区間を検出し、前記画像区間と音声区間それぞれを正規化することにより、入力された画像と音声とを同期させることを特徴とする請求項 7 に記載の学習装置。

【請求項 9】 前記同期手段は、さらに、正規化された前記画像区間と音声区間それぞれの始点と終点とを一致させることにより、入力された画像と音声とを同期させることを特徴とする請求項 8 に記載の学習装置。

【請求項 10】 前記同期手段は、さらに、正規化された前記画像区間と音声区間それぞれにおける画像のフレームと音声のフレームとを対応させることにより、入力された画像と音声とを同期させることを特徴とする請求項 8 に記載の学習装置。

【請求項 11】 入力された画像と音声とを同期させる同期ステップと、同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、前記抽出ステップにおいて出力される合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習ステップとを備えることを特徴とする学習方法。

【請求項 12】 所定の学習処理を、コンピュータに行わせるプログラムが記録されている記録媒体において、入力された画像と音声とを同期させる同期ステップと、同期された前記画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、前記抽出ステップにおいて出力される合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習ステップとを備えるプログラムが記録されていることを特徴とする記録媒体。

【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】本発明は、認識装置および認識方法、学習装置および学習方法、並びに記録媒体に関し、特に、例えば、画像と音声から、その概念を認識する場合において、高い認識性能を得ることができるようにする認識装置および認識方法、学習装置および学習方法、並びに記録媒体に関する。

## 【0002】

【従来の技術】近年においては、CPUの高速化、メモリの大容量化等が進み、例えば、音声認識装置を搭載した

エンタテインメント用のロボット等が低価格で実現されている。

【0003】このようなロボットは、ユーザの音声を音声認識し、その認識結果に基づいて、各種の行動を起こす。即ち、例えば、ユーザが、「お手」と発話した場合には、ロボットは、実際の犬がお手をするような行動を起こす。

## 【0004】

【発明が解決しようとする課題】ところで、例えば、ロボットにおいて、外部から与えられる刺激としての音声だけでなく、画像についても処理を行い、その音声と画像が表す概念を認識して、その認識した概念に基づいて、行動するようにすれば、よりエンタテインメント性の向上させたロボットを実現することができると予想される。

【0005】しかしながら、例えば、ロボットに対して、ある概念を表す音声と画像が与えられる場合、その音声と画像とは、同期したものとはなっていないため、そのような同期していない画像と音声を処理して、その画像と音声を表す概念を認識しても、十分な認識性能が得られないことが予想される。

【0006】また、上述のように、画像と音声を表す概念を認識するには、あらかじめ学習を行っておく必要があるが、学習に際しても、ある概念を表す音声と画像が、同期したものとはなっていない場合には、十分な認識性能が得られないことが予想される。

【0007】本発明は、このような状況に鑑みてなされたものであり、画像と音声から、その概念を認識する場合において、高い認識性能を得ることができるようにするものである。

## 【0008】

【課題を解決するための手段】本発明の認識装置は、入力された画像と音声を同期させる同期手段と、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出手段と、抽出手段において出力される合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声を表す概念を認識する認識手段とを備えることを特徴とする。

【0009】本発明の認識方法は、入力された画像と音

声を同期させる同期ステップと、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、抽出ステップにおいて出力される合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声を表す概念を認識する認識ステップとを備えることを特徴とする。

【0010】本発明の第1の記録媒体は、入力された画像と音声を同期させる同期ステップと、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、抽出ステップにおいて出力される合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声を表す概念を認識する認識ステップとを備えるプログラムが記録されていることを特徴とする。

【0011】本発明の学習装置は、入力された画像と音声を同期させる同期手段と、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出手段と、抽出手段において出力される合成特徴量と、合成特徴量に基づいて学習を行い、同一概念を表す画像および音声から得られる合成特徴量と、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習手段とを備えることを特徴とする。

【0012】本発明の学習方法は、入力された画像と音声を同期させる同期ステップと、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、抽出ステップにおいて出力される合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習ステップとを備えることを特徴とする。

【0013】本発明の第2の記録媒体は、入力された画像と音声を同期させる同期ステップと、同期された画像と音声それぞれから、特徴量を抽出し、その画像と音声の特徴量を合成した合成特徴量を出力する抽出ステップと、抽出ステップにおいて出力される合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する学習ステップとを備えるプログラムが記録されていることを特徴とする。

【0014】本発明の認識装置および認識方法、並びに第1の記録媒体においては、入力された画像と音声同期され、その同期された画像と音声それぞれから、特徴量が抽出され、その画像と音声の特徴量を合成した合成特徴量が出力される。そして、その合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことによ

り、入力された画像と音声を表す概念が認識される。

【0015】本発明の学習装置および学習方法、並びに第2の記録媒体においては、入力された画像と音声同期され、その同期された画像と音声それぞれから、特徴量が抽出され、その画像と音声の特徴量を合成した合成特徴量が出力される。そして、その合成特徴量に基づいて学習が行われることによりモデルが生成され、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書が生成される。

【0016】

【発明の実施の形態】図1は、本発明を適用した認識装置の一実施の形態の構成例を示している。

【0017】センサ部1は、例えば、少なくとも、音声(音)を集音して出力するマイク(マイクロフォン)と、画像を撮影して出力するビデオカメラから構成され、外部から与えられる刺激としての音声と画像を感知し、A/D(Analog Digital)変換して、対応する音声データと画像データを、同期処理部2に出力する。

【0018】同期処理部2は、センサ部1が出力する音声データと画像データについて、後述するような同期処理を施し、これにより、相互に同期した音声データと画像データを、所定のフレーム単位で、特徴抽出部3に供給する。

【0019】特徴抽出部3は、同期処理部2から供給される音声データを、そのフレーム(以下、適宜、音声フレームという)単位で処理し、音声の特徴パラメータを抽出する。さらに、特徴抽出部3は、同期処理部2から供給される画像データを、そのフレーム(以下、適宜、画像フレームという)単位で処理し、画像の特徴パラメータを抽出する。そして、特徴抽出部3は、音声フレームから得られた音声の特徴パラメータと、その音声フレームに対応する画像フレームから得られた画像の特徴パラメータとを合成し、その結果得られる合成特徴パラメータを、メモリ4に供給する。

【0020】なお、音声と画像の特徴パラメータの合成方法としては、例えば、その音声と画像の特徴パラメータがベクトルで構成される場合に、その音声と画像の特徴ベクトルの要素(コンポーネント)を接続して1つのベクトルを構成する方法等がある。

【0021】また、音声の特徴パラメータとしては、例えば、MFCC(Mel Frequency Cepstrum Coefficients)や、そのフレーム間差分、パワー等を用いることができる。さらに、画像の特徴パラメータとしては、動きベクトルや、色情報、DCT(Discrete Cosine Transform)係数、画像に表示された物体の形状を表す情報等を用いることができる。

【0022】メモリ4は、特徴抽出部3から供給される合成特徴パラメータを一時記憶する。

【0023】認識処理部5は、辞書データベース6の辞

10

20

30

40

50

書に登録されている各モデルと、メモリ4に記憶された合成特徴パラメータとを用いて、各モデルから、合成特徴パラメータが観測されるスコア(尤度)を求めるマッチング処理を行い、その合成特徴パラメータと最もマッチするモデル、即ち、例えば、最高スコアのモデルを求める。さらに、認識処理部5は、その最高スコアに基づいて、いま認識の対象となっている音声および画像に対応するモデルが、既に、辞書データベース6の辞書に登録されているかどうか、即ち、学習済みかどうかを判定し、学習済みの場合には、最高スコアのモデルに対応する概念情報を求め、その概念情報を、センサ部1に入力された音声と画像が表す概念の認識結果として出力する。また、認識処理部5は、いま認識の対象となっている音声および画像に対応するモデルが学習済みでない場合(いま認識の対象となっている音声および画像に対応するモデルが、辞書データベース6の辞書に登録されていない場合)、学習部7に対して、そのモデルの学習を要求する。

【0024】なお、ここでは、スコアは、その値が大きいほど尤度が高くなるものとする。また、スコアとしては、例えば、モデルから合成特徴パラメータが観測される確率や、合成特徴パラメータのモデルに対する類似性(例えば、特徴空間におけるモデルと合成特徴パラメータとの距離)(Confidence Measure)等を用いることが可能である。但し、スコアとして、例えば、特徴空間におけるモデルと合成特徴パラメータとの距離を用いる場合には、スコアの値が小さいほど、尤度が高いことを意味することになる。

【0025】辞書データベース6は、合成特徴パラメータを用いて学習を行うことにより得られる音声および画像のモデルと、その音声および画像の概念を表す概念情報とが対応付けられて登録される辞書を記憶している。なお、音声および画像のモデルとしては、例えば、HMM(Hidden Markov Model)等の確率モデルその他を採用することができる。

【0026】学習部7は、認識処理部5から学習の要求があると、メモリ4に記憶された合成特徴パラメータを用いて、音声および画像のモデルの学習を行う。さらに、学習部7は、学習の結果得られたモデルを、そのモデルが表す音声および画像の概念情報と対応付けて、辞書データベース6の辞書に登録する。

【0027】入力部8は、例えば、キーボード等で構成され、学習部7が学習の対象としているモデルに対応する音声および画像の概念情報を入力するときに、ユーザによって操作される。入力部8が操作されることにより入力された概念情報は、学習部7に供給される。なお、入力部8は、その他、例えば、音声認識装置等で構成することも可能であり、この場合、概念情報は、音声で入力することができる。

【0028】次に、図2のフローチャートを参照して、

図1の認識装置の動作について説明する。

【0029】センサ部1は、音声と画像を感知し、対応する音声データと画像データを、同期処理部2に出力する。同期処理部2は、ステップS1において、センサ部1が出力する音声データと画像データについて、同期処理を施し、これにより、相互に同期した、所定の音声フレームごとの音声データと、所定の画像フレームごとの画像データを、特徴抽出部3に供給する。

【0030】特徴抽出部3は、ステップS2において、同期処理部2からの音声フレーム単位の音声データから、その特徴パラメータを抽出するとともに、同じく、同期処理部2からの画像フレーム単位の画像データから、その特徴パラメータを抽出する。さらに、特徴抽出部3は、ステップS2において、各音声フレームから得られた音声の特徴パラメータと、その音声フレームに対応する画像フレームから得られた画像の特徴パラメータとを合成し、合成特徴パラメータとして、メモリ4に供給して記憶させる。

【0031】そして、ステップS3に進み、認識処理部5は、マッチング処理を行い、スコアを求める。即ち、認識処理部5は、辞書データベース6の辞書に登録されている各モデルと、メモリ4に記憶された合成特徴パラメータとを用いて、各モデルから、合成特徴パラメータが観測されるスコアを求め、さらに、各モデルのスコアのうち、最も値の大きいもの（最高スコア）を求めて、ステップS4に進む。

【0032】ステップS4では、認識処理部5は、最高スコアが所定の閾値 $\varepsilon$ よりも大（所定の閾値以上）であるかどうかによって、センサ部1に入力された音声および画像のモデルが学習済みであるかどうかを判定する。ここで、モデルが学習済みかどうかの判定は、その他、例えば、音声認識において、入力音声が、辞書に登録されていない単語（未知語（OOV (Out Of Vocabulary)））かどうかを判定する手法等を適用して行うことが可能である。

【0033】ステップS4において、最高スコアが所定の閾値 $\varepsilon$ よりも大であると判定された場合、認識処理部5は、センサ部1に入力された音声および画像のモデルが学習済みであるとして、ステップS5に進む。ステップS5では、認識処理部5は、最高スコアのモデルに対応付けられた概念情報を、辞書データベース6の辞書から読み出し、認識結果として出力して、処理を終了する。

【0034】一方、ステップS4において、最高スコアが所定の閾値 $\varepsilon$ よりも大でないと判定された場合、認識処理部5は、いま認識の対象となっている音声および画像のモデルが学習済みでないと、学習部7に対して、そのモデルの学習を要求して、ステップS6に進む。

【0035】ステップS6では、学習部7は、メモリ4

に記憶された合成特徴パラメータを用いて、センサ部1に入力された音声および画像のモデルの学習を行う。なお、モデルとしてHMMを採用する場合には、モデルの学習には、例えば、Baum-Welchの再推定法などを用いることが可能である。

【0036】学習部7は、学習によってモデルを得ると、そのモデルが表す音声および画像の概念情報の入力を要求するメッセージを、図示せぬディスプレイまたはスピーカより出力し、ユーザが、入力部8を操作することにより、概念情報を入力するのを待って、ステップS7に進む。

【0037】ステップS7では、学習部7は、学習の結果得られたモデルと、入力部8が操作されることにより入力された概念情報と対応付けて、辞書データベース6の辞書に追加登録し、処理を終了する。

【0038】なお、上述の場合には、最高スコアと閾値 $\varepsilon$ との大小関係に基づいて、認識処理部5から認識結果を出力し、あるいは、学習部7において学習を行うようにしたが、認識処理部5から認識結果を出力するか、あるいは、学習部7において学習を行うかは、その他、例えば、装置において、認識モードと学習モードとを切り替えるスイッチ等を設け、そのスイッチの操作に基づいて決定することが可能である。

【0039】次に、図2のステップS1において同期処理部2が行う同期処理について説明する。

【0040】例えば、いま、「ボールを蹴る」という音声と、ボールを蹴っている状態の画像とを、センサ部1から入力するとした場合、「ボールを蹴る」という音声が存在する区間（音声区間）と、ボールを蹴っている状態の画像（ボールを蹴るという行動が行われている画像）が存在する区間（画像区間）とは、図3に示すように、同期した状態にない。

【0041】即ち、「ボールを蹴る」という音声の音声区間の始点（時刻） $S_b$ と、ボールを蹴っている状態の画像の画像区間の始点 $M_b$ とは、一般に一致せず、「ボールを蹴る」という音声の音声区間の終点（時刻） $S_e$ と、ボールを蹴っている状態の画像の画像区間の終点 $M_e$ も、一般に一致しない。従って、その音声区間の長さ $T_s (= S_e - S_b)$ と、画像区間の長さ $T_v (= M_e - M_b)$ も一致しない。

【0042】また、上述のように、認識モードと学習モードとをスイッチによって切り替える場合に、学習モードにおいて学習を行おうとするときには、その学習対象の音声および画像は、例えば、図4に示すように、繰り返し入力されることがある。この場合も、図3における場合と同様に、繰り返し入力される音声の音声区間と、画像の画像区間とは、同期した状態にはならない。

【0043】即ち、いま、 $i$ 番目に入力される音声と画像を、それぞれ $S_i$ と $M_i$ と表すとともに、音声 $S_i$ の音声区間の始点、終点、長さを、それぞれ $B(S_i)$ 、 $E$

10

20

30

40

50

( $S_i$ ),  $T(S_i)$ と表し、画像 $M_i$ の画像区間の始点、終点、長さを、それぞれ $B(M_i)$ ,  $E(M_i)$ ,  $T(M_i)$ と表すと、 $i$ 番目の音声区間の始点 $B(S_i)$ と画像区間の始点 $B(M_i)$ は、一般に一致せず、終点 $E(S_i)$ と $E(M_i)$ も、一般には一致しない。従って、音声区間の長さ $T(S_i)$ と、画像区間の長さ $T(M_i)$ も、一般には一致しない。

【0044】このように、音声区間と画像区間の始点、終点、長さが一致しないと、音声と画像の特徴パラメータを合成した合成特徴パラメータにおいて、音声または画像の特徴パラメータのうちのいずれか一方が存在しない(いずれか一方だけ存在する)区間が生じ、この場合、ある同一の概念について、音声および画像の両方が入力されているのに、その概念の認識に、音声または画像のうちの一方だけしか用いられず、その結果、音声および画像の両方が用いられる場合に比較して、認識性能が劣化することとなる。

【0045】そこで、図1の同期処理部2は、図5に示すように、センサ部1から供給される音声 $S$ と画像 $M$ を同期させる。

【0046】即ち、同期処理部2は、例えば、図5(A)に示すように、音声 $S$ の音声区間の始点 $S_b^*$ と、画像 $M$ の画像区間の始点 $M_b^*$ とを、ある時刻に一致させるとともに、音声区間の終点 $S_e^*$ と、画像区間の終点 $M_e^*$ とを、他の時刻に一致させる正規化処理を行う。なお、正規化が行われることにより、音声区間の長さ $T(S_i)$ と、画像区間の長さ $T(M_i)$ とは、一致することになる。

【0047】より具体的には、同期処理部2は、例えば、音声の音声区間の始点 $S_b^*$ 、または画像 $M$ の画像区間の始点 $M_b^*$ のうちのいずれか一方を基準点とし、その基準点に、音声区間の始点 $S_b^*$ と、画像区間の始点 $M_b^*$ を一致させる。さらに、同期処理部2は、例えば、200、400、または800ミリ秒等の所定の時間を基準時間長とし、音声区間と画像区間の長さが、いずれも基準時間長に一致するように、音声区間の終点 $S_e^*$ と、画像区間の終点 $M_e^*$ を変更する。従って、音声区間の終点 $S_e^*$ と、画像区間の終点 $M_e^*$ は、基準点から基準時間長だけ経過した点において一致することになる。

【0048】さらに、同期処理部2は、正規化された音声区間の各音声フレームと、画像区間の各画像フレームとを、例えば、図5(B)に示すように、一対一に対応させる対応付け処理を行い、これにより、音声と画像とを同期させる。

【0049】なお、正規化処理によれば、音声区間または画像区間の長さが変更することから、その音声区間または画像区間を構成する音声フレームまたは画像フレームの数を増減させる必要がある。また、対応付け処理においても、音声フレームと画像フレームの時間長が異なる場合には、音声フレームまたは画像フレームの数を増減させる必要があることがある。この音声フレームまた

は画像フレームの数の増減は、例えば、補間や間引き等によって行うことが可能である。

【0050】ここで、上述の場合には、音声フレームと画像フレームとを、一対一に対応付けるようにしたが、音声フレームと画像フレームとは、一対多または多対一に対応付けることも可能である。即ち、例えば、画像フレームの時間長が、音声フレームの時間長の $L$ 倍に一致する場合には、1フレームの画像フレームと、 $L$ フレームの音声フレームとを対応付けるようにすることが可能である。

【0051】次に、図6は、図1の同期処理部2の構成例を示している。

【0052】センサ部1(図1)が出力する画像と音声は、区間検出部11とメモリ12に供給される。

【0053】区間検出部11は、そこに供給される画像の画像区間と、音声の音声区間とを検出し、区間正規化部13および同期化部14に供給する。即ち、区間検出部11は、例えば、各画像フレーム全体の動きベクトルを求め、その動きベクトルに基づき、ある程度の動きのある画像フレームが連続している区間を画像区間として検出する。また、区間検出部11は、例えば、各音声フレームのパワーを求め、そのパワーに基づき、ある程度のパワーを有する音声フレームが連続している区間を音声区間として検出する。

【0054】メモリ12は、そこに供給される画像データと音声データを一時記憶する。

【0055】区間正規化部13は、区間検出部11から供給される音声区間を構成する音声フレームの音声データと、同じく区間検出部11から供給される画像区間を構成する画像フレームの画像データを、メモリ12から読み出す。さらに、区間正規化部13は、その音声区間の音声データと、画像区間の画像データについて、上述した正規化処理を行い、これにより、始点、終点、長さを一致させた(正規化された)音声区間と画像区間の音声データと画像データを得て、同期化部14に供給する。

【0056】同期化部14は、区間正規化部13からの音声データと画像データについて、上述の対応付け処理を行うことにより、正規化された音声区間と画像区間の音声フレームと画像フレームとを一対一に対応付け、特徴抽出部3(図1)に出力する。

【0057】次に、図7のフローチャートを参照して、図6の同期処理部2が、図2のステップS1において行う同期処理について説明する。

【0058】センサ部1(図1)が出力する音声データと画像データは、区間検出部11とメモリ12に供給され、メモリ12は、その音声データと画像データを一時記憶する。

【0059】区間検出部11は、ステップS11において、センサ部1からの画像データの画像区間と、音声デ

10

20

30

40

50

ータの音声区間とを検出し、区間正規化部13および同期化部14に供給して、ステップS12に進む。

【0060】ステップS12では、区間正規化部13は、区間検出部11からの音声区間と画像区間それぞれを構成する音声フレームの音声データと、画像フレームの画像データを、メモリ12から読み出し、正規化処理を施すことで、始点、終点、長さを一致させた音声区間と画像区間、即ち、正規化された音声区間と画像区間の音声データと画像データを得て、同期化部14に供給する。

【0061】同期化部14は、ステップS13において、区間正規化部13からの、正規化された音声区間と画像区間の音声フレームと画像フレームとを一对一に対応付ける対応付け処理を行うことで、音声フレームと画像フレームとをフレーム同期化し、特徴抽出部3(図1)に出力して、同期処理を終了する。

【0062】以上のように、センサ部1から入力される、同一概念を表す音声データと画像データとを同期させるようにしたので、その音声データと画像データから得られる特徴パラメータを合成した合成特徴パラメータにおいて、音声または画像の特徴パラメータのうちのいずれか一方が存在しない区間が生じなくなり、そのような合成特徴パラメータを用いて認識が行われる結果、認識性能を向上させることができる。

【0063】次に、上述した一連の処理は、ハードウェアにより行うこともできるし、ソフトウェアにより行う場合にも、そのソフトウェアを構成するプログラムが、汎用のコンピュータ等にインストールされる。

【0064】そこで、図8は、上述した一連の処理を実行するプログラムがインストールされるコンピュータの一実施の形態の構成例を示している。

【0065】プログラムは、コンピュータに内蔵されている記録媒体としてのハードディスク105やROM103に予め記録しておくことができる。

【0066】あるいはまた、プログラムは、フロッピー(登録商標)ディスク、CD-ROM(Compact Disc Read Only Memory)、MO(Magneto optical)ディスク、DVD(Digital Versatile Disc)、磁気ディスク、半導体メモリなどのリムーバブル記録媒体111に、一時的あるいは永続的に格納(記録)しておくことができる。このようなリムーバブル記録媒体111は、いわゆるパッケージソフトウェアとして提供することができる。

【0067】なお、プログラムは、上述したようなリムーバブル記録媒体111からコンピュータにインストールする他、ダウンロードサイトから、デジタル衛星放送用の人工衛星を介して、コンピュータに無線で転送したり、LAN(Local Area Network)、インターネットといったネットワークを介して、コンピュータに有線で転送し、コンピュータでは、そのようにして転送されてくる

プログラムを、通信部108で受信し、内蔵するハードディスク105にインストールすることができる。

【0068】コンピュータは、CPU(Central Processing Unit)102を内蔵している。CPU102には、バス101を介して、入出力インタフェース110が接続されており、CPU102は、入出力インタフェース110を介して、ユーザによって、キーボードや、マウス、マイク等で構成される入力部107が操作等されることにより指令が入力されると、それにしたがって、ROM(Read Only Memory)103に格納されているプログラムを実行する。あるいは、また、CPU102は、ハードディスク105に格納されているプログラム、衛星若しくはネットワークから転送され、通信部108で受信されてハードディスク105にインストールされたプログラム、またはドライブ109に装着されたリムーバブル記録媒体111から読み出されてハードディスク105にインストールされたプログラムを、RAM(Random Access Memory)104にロードして実行する。これにより、CPU102は、上述したフローチャートにしたがった処理、あるいは上述したブロック図の構成により行われる処理を行う。そして、CPU102は、その処理結果を、必要に応じて、例えば、入出力インタフェース110を介して、LCD(Liquid Crystal Display)やスピーカ等で構成される出力部106から出力、あるいは、通信部108から送信、さらには、ハードディスク105に記録等させる。

【0069】ここで、本明細書において、コンピュータに各種の処理を行わせるためのプログラムを記述する処理ステップは、必ずしもフローチャートとして記載された順序に沿って時系列に処理する必要はなく、並列的あるいは個別に実行される処理(例えば、並列処理あるいはオブジェクトによる処理)も含むものである。

【0070】また、プログラムは、1のコンピュータにより処理されるものであっても良いし、複数のコンピュータによって分散処理されるものであっても良い。さらに、プログラムは、遠方のコンピュータに転送されて実行されるものであっても良い。

【0071】なお、図1の認識装置は、例えば、エンタテインメント用のロボットや、音声対話システム等に適用可能である。図1の認識装置をロボットに適用した場合には、例えば、ユーザの発話(例えば、「ボールが転がっている」など)と、対応する画像(例えば、ボールが転がっている様子が撮影された画像)から、周囲で生じている事象を、より正確に表す概念を得て、ロボットに的確な行動をとらせること等が可能となる。また、図1の認識装置を対話システムに適用した場合には、例えば、ユーザの発話とジェスチャから、ユーザの意図を、より正確に表す概念を得て、的確な応答を返すようにすること等が可能となる。

【0072】ここで、本実施の形態においては、音声と



画像の両方から、それらが表す概念を認識するようにしたが、音声または画像のいずれか一方だけから、それが表す概念を認識するようにすることも可能である。

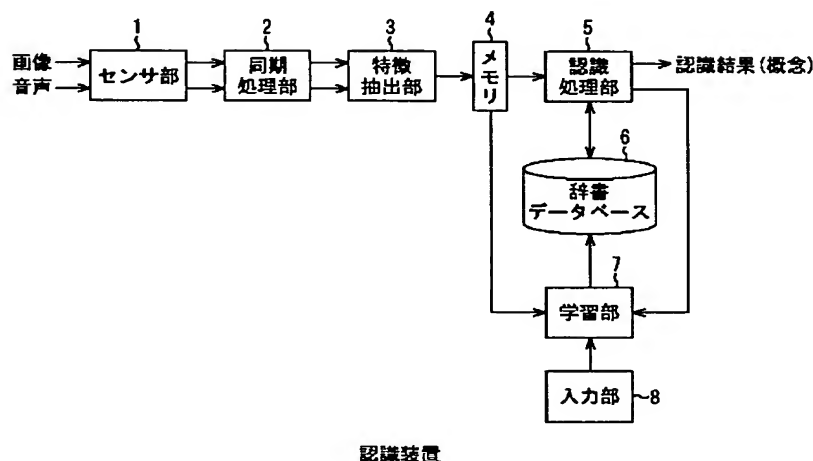
【0073】また、センサ部1には、マイクとビデオカメラだけでなく、例えば、圧力センサ等の、画像および音声以外の外部からの刺激を感知するセンサを設け、そのセンサで得られる情報も、合成特徴パラメータに含めて、認識および学習を行うようにすることが可能である。

【0074】

【発明の効果】本発明の認識装置および認識方法、並びに第1の記録媒体によれば、入力された画像と音声を同期させ、その同期された画像と音声それぞれから、特徴量を抽出して、その画像と音声の特徴量を合成した合成特徴量を得る。そして、その合成特徴量と、辞書におけるモデルとを用いてマッチングを行うことにより、入力された画像と音声を表す概念を認識する。従って、認識性能を向上させることが可能となる。

【0075】本発明の学習装置および学習方法、並びに第2の記録媒体によれば、入力された画像と音声を同期させ、その同期された画像と音声それぞれから、特徴量を抽出して、その画像と音声の特徴量を合成した合成特徴量を得る。そして、その合成特徴量に基づいて学習を行うことによりモデルを生成し、同一概念を表す画像および音声に対応するモデルと、その画像および音声の概念を表す概念情報とを対応付けた辞書を生成する。従って、その辞書を用いた概念の認識を行う場合において、その認識性能を向上させることが可能となる。

【図1】



【図面の簡単な説明】

【図1】本発明を適用した認識装置の一実施の形態の構成例を示すブロック図である。

【図2】認識装置の処理を説明するフローチャートである。

【図3】音声と画像が同期していない様子を示す図である。

【図4】音声と画像が同期していない様子を示す図である。

【図5】音声と画像が同期している様子を示す図である。

【図6】同期処理部2の構成例を示すブロック図である。

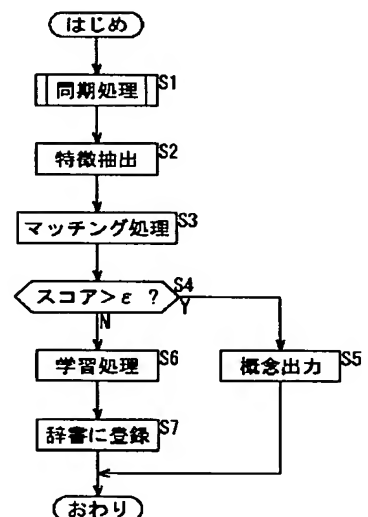
【図7】同期処理部2による同期処理を説明するフローチャートである。

【図8】本発明を適用したコンピュータの一実施の形態の構成例を示すブロック図である。

【符号の説明】

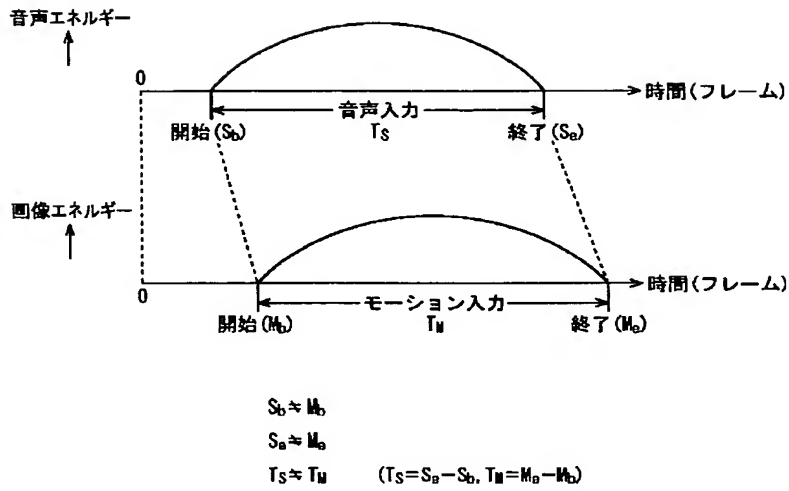
1 センサ部, 2 同期処理部, 3 特徴抽出部, 4 メモリ, 5 認識処理部, 6 辞書データベース, 7 学習部, 8 入力部, 11 区間検出部, 12 メモリ, 13 区間正規化部, 14 同期化部, 101 バス, 102 CPU, 103 ROM, 104 RAM, 105 ハードディスク, 106 出力部, 107 入力部, 108 通信部, 109 ドライブ, 110 入出力インタフェース, 111 リムーバブル記録媒体

【図2】

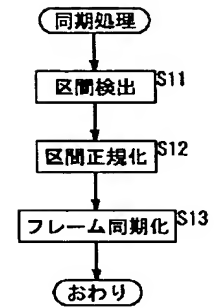




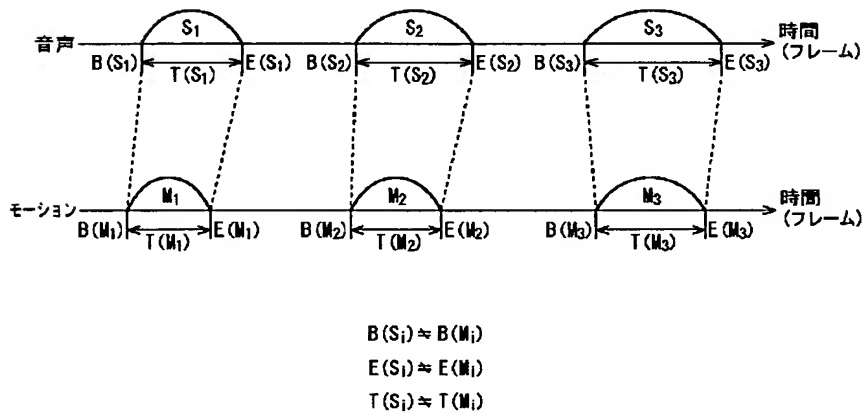
【図3】



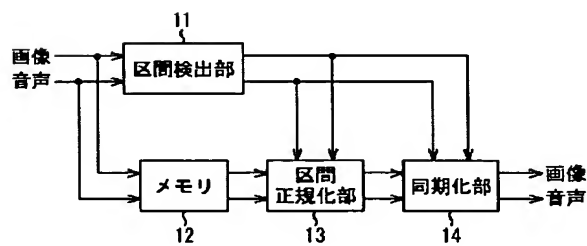
【図7】



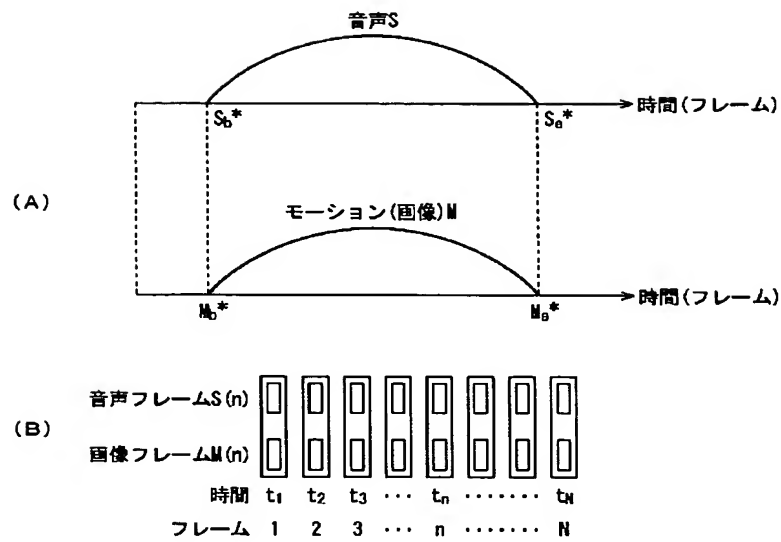
【図4】



【図6】



【図5】



【図8】

